

# Architectural Optimization of the Monad Layer 1 Blockchain

Nikhil Ranjan  
nikhil@polkassembly.io

The evolution of decentralized smart contract platforms has reached a critical bottleneck where the historical design choices of the Ethereum Virtual Machine (EVM) severely limit transaction execution speed.<sup>1</sup> The sequential nature of the EVM execution client prevents modern multi-core processors from being fully utilized, capping execution throughput and stretching latency tails under network congestion.<sup>1</sup> Monad has emerged as an independent Layer 1 blockchain seeking to break this transaction throughput ceiling by delivering up to 10,000 transactions per second (TPS) with a 400 ms block frequency and an 800 ms finality time, while maintaining 100% bytecode compatibility with Ethereum.<sup>4</sup>

Founded by former Jump Trading engineers Keone Hon and James Hunsaker, Monad applies high-frequency trading (HFT) optimization principles to blockchain architecture.<sup>6</sup> This high performance is achieved by rearchitecting the execution, storage, and consensus pipelines through five major innovations: a custom MonadBFT consensus engine, a RaptorCast block propagation protocol, asynchronous execution, optimistic parallel transaction processing, and a native Patricia trie state database known as MonadDB.<sup>3</sup> However, pushing physical hardware and distributed network limits to these extremes introduces deep systemic trade-offs, physical performance cliffs, and unique engineering vulnerabilities. This analysis provides an exhaustive technical critique of Monad's architecture, evaluates its performance under high load, and identifies key pathways for technological optimization.

---

## Consensus and Block Propagation Optimization: MonadBFT and RaptorCast

At the foundation of Monad's high-throughput architecture is MonadBFT, a Byzantine Fault Tolerant (BFT) consensus protocol derived from the HotStuff family.<sup>4</sup> Traditional BFT consensus algorithms, such as Practical Byzantine Fault Tolerance (PBFT) and Tendermint, exhibit quadratic communication complexity ( $O(n^2)$ ) due to all-to-all communication topologies.<sup>9</sup> This design severely restricts the size of the active validator set, as messaging overhead becomes unsustainable when scaling beyond a few dozen nodes.<sup>4</sup> MonadBFT solves this scalability constraint by achieving linear message and authenticator complexity ( $O(n)$ ) on the

happy path, enabling the protocol to support a decentralized network of 150 to 200 active validators on modest hardware.<sup>4</sup>

The protocol achieves this linear complexity by utilizing an optimistic, pipelined "one-to-many-to-one" communication structure.<sup>6</sup> In each consensus round (or view), a preselected leader proposes a block and collects cryptographic votes (attestations) from the validator set.<sup>6</sup> Rather than waiting for multiple round-trips to finalize a block before proposing the next one, MonadBFT pipelines the commit stages.<sup>4</sup> A block enters the consensus pipeline in each round; while block  $N$  is being finalized over successive rounds, the leader is already proposing and collecting votes for block  $N + 1$ .<sup>4</sup> This overlapping progression allows the protocol to operate at network latency, achieving speculative finality within a single round and complete cryptographic finality within two rounds (~800 ms).<sup>4</sup>

A significant architectural contribution of MonadBFT is its built-in resistance to "tail-forking".<sup>8</sup> In standard pipelined BFT protocols, if a leader goes offline or fails to collect a supermajority of votes, the previous unfinalized proposal is often abandoned.<sup>8</sup> Malicious leaders can exploit this vulnerability as a form of Maximal Extractable Value (MEV) attack, deliberately discarding their predecessor's block.<sup>8</sup> This deprives the previous proposer of rewards and allows the malicious leader to reorder, censor, or insert transactions across blocks to capture arbitrage opportunities.<sup>8</sup>

MonadBFT mitigates tail-forking by enforcing a strict block reproposal mechanism during view changes (the unhappy path).<sup>8</sup> When a consensus round fails due to network delays or an offline leader, validators include their "local tip" (the latest unfinalized proposal they voted for) in their timeout messages.<sup>8</sup> The next leader is mathematically obligated to aggregate these local tips, identify the "high tip" (the block from the highest view), and repropose it.<sup>8</sup> If the leader cannot recover the high-tip block, validators jointly issue a cryptographic No-Endorsement Certificate (NEC), which guarantees the block could never have achieved a Quorum Certificate (QC), allowing the leader to propose a new block safely.<sup>13</sup>

This high-frequency block generation creates a massive data propagation hurdle.<sup>4</sup> At 10,000 TPS and an average transaction size of 200 bytes, a validator must transmit 2 MB of block data per second.<sup>4</sup> Directly uploading a 2 MB block to 200 validators would require a leader to have 400 MB/s (3.2 Gbps) of upload bandwidth, which exceeds the capacity of standard validator hardware.<sup>4</sup> Monad resolves this bandwidth bottleneck through RaptorCast, a specialized block propagation protocol.<sup>4</sup>

RaptorCast rejects TCP in favor of UDP to minimize connection overhead and lower transmission latency.<sup>14</sup> It employs erasure coding based on Raptor codes, where a block is split into a series of smaller chunks.<sup>4</sup> The total size of all chunks is greater than the original block, but the original block can be fully reconstructed from any subset of chunks whose aggregate size matches the block's original volume (e.g., reclaiming a 1000 KB block from any 50 out of 150

generated chunks).<sup>10</sup> RaptorCast distributes these chunks using a structured, two-level broadcast tree.<sup>4</sup> The leader distributes unique chunks to "level-1" validators (assigned proportionally to their stake weight), who then rebroadcast their assigned chunk to all "level-2" validators.<sup>6</sup> This design limits worst-case propagation time to twice the worst-case one-way network latency while keeping the leader's upload bandwidth independent of the total validator count.<sup>14</sup>

Parameter	Traditional PBFT / Tendermint	MonadBFT / RaptorCast
Communication Complexity	Quadratic ( $O(n^2)$ ) <sup>9</sup>	Linear ( $O(n)$ ) on happy path <sup>4</sup>
Active Validator Capacity	Low ( $\leq$ nodes for acceptable latency) <sup>9</sup>	High (150–200 nodes supported) <sup>4</sup>
Commit Latency	Sequential, multi-round block intervals <sup>15</sup>	Pipelined, single-round speculative, two-round final <sup>4</sup>
Block Time / Finality	12 seconds (Ethereum) <sup>6</sup>	400 ms block / 800 ms finality <sup>4</sup>
MEV Tail-Forking Resistance	Vulnerable (Unfinalized blocks easily discarded) <sup>8</sup>	High (Enforced reproposal of high tip via TCs/NECs) <sup>8</sup>
Transmission Protocol	TCP gossip networks (High latency hops) <sup>14</sup>	UDP Raptor code erasure chunking <sup>10</sup>
Required Leader Upload Bandwidth	Scaled linearly with validator count ( $O(n)$ )	Scaled independently of validator count via 2-level tree <sup>10</sup>

---

## The Deferred Execution Pipeline and the RPC Consistency Paradox

A core structural limitation of legacy EVM blockchains is the coupling of consensus and execution.<sup>12</sup> In Ethereum, the block proposer must execute all transactions in a proposed block before sharing the proposal, and validating nodes must execute them before responding with

their vote.<sup>15</sup> This design squeezes the available execution window, forcing the system to allocate most of the block time to network communication and safety margins rather than computing transactions.<sup>6</sup> Monad decouples consensus from execution through an asynchronous execution model.<sup>6</sup>

In Monad, validators reach consensus on the linear ordering of transactions in a block *before* executing them.<sup>5</sup> Once transaction ordering is finalized via MonadBFT, execution occurs asynchronously on a separate pipeline, allowing execution to utilize the full block time without blocking consensus.<sup>4</sup> While execution workers process block  $N$ , consensus can proceed with voting on and proposing block  $N + 1$ .<sup>3</sup> This decoupling allows Monad to increase its block gas limit to 150 million gas in testnet (equivalent to a gas rate of 375 million gas/second) and support highly complex smart contracts.<sup>4</sup>

However, this asynchronous execution model introduces a severe systems engineering challenge: the **loss of immediate read-your-own-writes (RYW) consistency** at the JSON-RPC layer.<sup>17</sup> In a standard EVM node, the state database is updated in lockstep with consensus, ensuring that client reads and writes are perfectly synchronized.<sup>15</sup> Because Monad defers execution, the state database on standard nodes lags behind consensus.<sup>3</sup> This replication lag breaks immediate RYW consistency, which is a core expectation for client-side applications, trading bots, and keepers.<sup>17</sup>

When a user submits a transaction via `eth_sendRawTransaction`, the RPC node cannot immediately validate nonces or balances because execution is lagging.<sup>17</sup> Consequently, invalid transactions (e.g., those with insufficient gas or incorrect nonces) are accepted into the mempool and only rejected much later during execution, complicating client-side error handling.<sup>17</sup> Furthermore, standard Geth RPC queries like `eth_getTransactionByHash` return null for pending transactions, as the transaction lifecycle behaves differently under the deferred execution model.<sup>17</sup>

If a trading bot submits a swap transaction and immediately queries the state database to read its new balance or update its pricing models, the query may route to a read replica or an execution thread that is still processing the previous block's state delta.<sup>17</sup> This lag breaks read-after-write verification, leading to stale reads and false-negative balances.<sup>18</sup>

This execution lag also impacts historical data access and event indexing.<sup>17</sup> Because Monad's high throughput generates up to 6TB of historical data per day, standard full nodes aggressively prune state, retaining only a rolling window of approximately 40,000 blocks (~4.5 hours of history).<sup>17</sup>

Consequently, deep historical `eth_call` and state queries are unavailable on standard nodes.<sup>17</sup> Furthermore, `eth_getLogs` requests are severely throttled to tight block ranges (typically 100 to 1,000 blocks per request) to prevent concurrent queries from overwhelming the node's memory.<sup>17</sup> This forces developers to rely on specialized, sharded archive architectures (such as

Goldsky Edge) to index and serve historical data, introducing external dependencies into the developer workflow.<sup>20</sup>

RPC API Feature	Ethereum (Synchronous EVM)	Monad (Asynchronous Deferred EVM)
Execution State Sync	Strictly synchronous (State updates with consensus) <sup>15</sup>	Eventually consistent (State lags behind consensus) <sup>3</sup>
Transaction Submission	Real-time state checks (Nonces and balances validated instantly)	Speculative state checks (Validation occurs post-consensus) <sup>10</sup>
eth_getTransactionByHash	Returns transaction details and status immediately	May return null for pending transactions during execution lag <sup>17</sup>
Local Node State Depth	Full state history retained by default	Aggressively pruned (Only ~40,000 blocks retained) <sup>17</sup>
eth_getLogs Range Limit	High (Often up to 10,000+ blocks per request)	Heavily restricted (100–1,000 blocks per request) <sup>17</sup>
Read-Your-Own-Writes	Guaranteed on the same node	Broken without session pinning or sticky routing <sup>18</sup>
Tooling Dependency	Low (Standard Geth RPC is sufficient)	High (Requires specialized indexers like Goldsky Edge) <sup>20</sup>

## Dynamic Parallel Execution and Runtime Re-Execution Overheads

To maximize CPU utilization, Monad implements an Optimistic Concurrency Control (OCC) model for its parallel execution engine.<sup>21</sup> Transactions in a block are assumed to be independent and are processed in parallel across multiple CPU cores.<sup>2</sup> This dynamic approach stands in sharp contrast to static parallel execution models, such as Solana’s Sealevel and Sui’s object model, which require transactions to explicitly declare all read and write dependencies

upfront.<sup>2</sup>

While Solana and Sui avoid runtime conflicts by scheduling independent transactions in advance, they force developers to abandon standard EVM designs and rewrite contracts to handle explicit state declarations.<sup>1</sup> Monad preserves the linear transaction ordering of Ethereum.<sup>4</sup> It dynamically parallelizes execution without requiring any code changes from developers.<sup>2</sup>

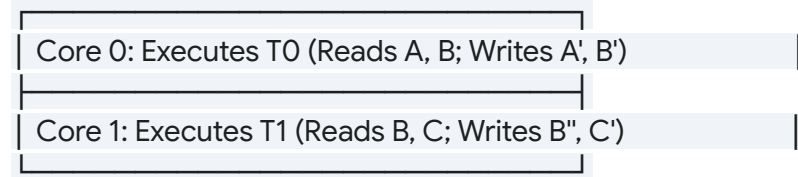
During parallel execution (Phase 1), Monad's runtime schedules transactions across separate threads, recording all read operations (inputs) and write operations (outputs) in a temporary memory structure called a PendingResult.<sup>25</sup> During the serial commit phase (Phase 2), transaction outputs are sequentially merged into the global state.<sup>21</sup> Before applying a transaction's updates, the system checks whether its inputs match the current state.<sup>25</sup> If no prior transaction in the sequence has modified the state keys read by this transaction, the outputs are committed.<sup>25</sup>

If a conflict is detected (i.e., a prior transaction has modified a state key that a subsequent transaction read), the subsequent transaction's execution is aborted and re-executed with the updated state.<sup>21</sup> This dynamic dependency management ensures that the final state matches the result of sequential execution.<sup>5</sup>

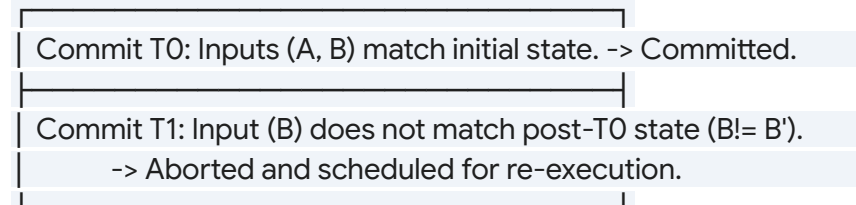
Transaction 0: Transfer A -> B (No dependencies)

Transaction 1: Transfer B -> C (Depends on balance of B updated by Transaction 0)

Phase 1: Optimistic Parallel Run (T0 and T1 execute simultaneously)



Phase 2: Serial Commitment and Conflict Validation



While Monad minimizes re-execution overhead by caching state in memory and avoiding

re-running non-state-dependent tasks (such as signature recovery)<sup>21</sup>, dynamic tracking still introduces significant runtime computation.<sup>22</sup>

In highly active DeFi environments with frequent shared state modifications (such as liquidity pools, shared AMM states, or liquidations), state hotspots are common.<sup>26</sup> When multiple transactions target the same state keys, Monad's parallel execution model can suffer from high abort rates.<sup>26</sup>

Under heavy state contention, the sequential validation phase is forced to discard the optimistic parallel runs and execute transactions sequentially.<sup>26</sup> This sequential bottleneck can drag down performance, potentially canceling out the benefits of parallel execution.<sup>26</sup>

To analyze how Monad compares to other high-performance blockchains, the table below outlines the trade-offs of different execution models:

<b>Dimension</b>	<b>Solana (Sealevel)</b>	<b>Sui (Move VM)</b>	<b>SupraBTM (iBTM)</b>	<b>Monad (Parallel EVM)</b>
<b>Dependency Management</b>	Static (Pre-declared in transaction signatures) <sup>2</sup>	Static (Object ownership metadata) <sup>23</sup>	Static deployment analysis & pre-run DAG <sup>26</sup>	Dynamic (Dynamic tracking at runtime) <sup>24</sup>
<b>State Paradigm</b>	Account-based (Global state) <sup>23</sup>	Object-centric (Owned/Shared objects) <sup>23</sup>	Account-based (EVM-native) <sup>26</sup>	Account-based (EVM-equivalent) <sup>5</sup>
<b>Conflict Resolution</b>	Prevented upfront by block scheduler <sup>24</sup>	Prevented upfront (No shared-state locks) <sup>23</sup>	Eliminated upfront via DAG scheduling <sup>26</sup>	Runtime abort, rollback, and sequential retry <sup>22</sup>
<b>Re-Execution Risk</b>	None <sup>24</sup>	None <sup>23</sup>	Minimal <sup>26</sup>	High under high state contention <sup>26</sup>
<b>High Contention Performance</b>	Stable (Bounded by scheduler queues)	High (Optimized for independent transactions) <sup>23</sup>	High (Dynamic fallback to sequential) <sup>26</sup>	Degrades to sequential speed with high compute overhead <sup>26</sup>

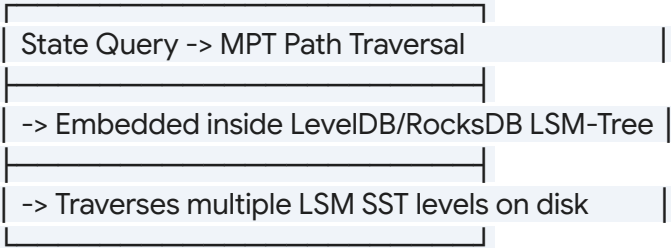
<b>Developer UX</b>	High complexity (Must declare all accounts upfront) <sup>2</sup>	High complexity (Must adapt to object ownership rules) <sup>28</sup>	Low complexity (Standard EVM compatibility)	Low complexity (100% bytecode and tooling equivalence) <sup>3</sup>
---------------------	--	--	---	---

## MonadDB: Custom Radix Tree and Filesystem Bypass

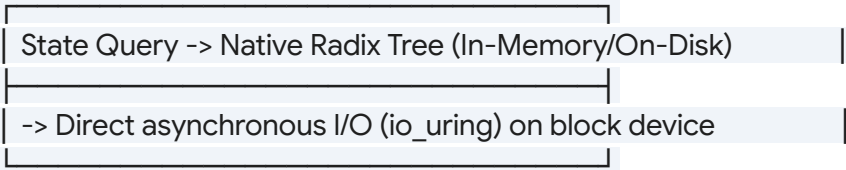
Execution speed is wasted if storage cannot keep up.<sup>1</sup> Most traditional Ethereum clients use generic key-value databases, such as LevelDB or RocksDB (implemented as LSM-Trees) or LMDB (implemented as B-Trees), to store state.<sup>30</sup> However, Ethereum structures its state as a Merkle Patricia Trie (MPT) to generate cryptographic proofs.<sup>30</sup>

This creates a suboptimal architecture where the MPT data structure is embedded inside a B-Tree or LSM-Tree database.<sup>30</sup> To read a single state value, the database must perform nested disk reads across multiple levels of the LSM-Tree or B-Tree, creating a severe I/O bottleneck.<sup>4</sup>

### Traditional EVM Storage (Double Layer Indirection):



### MonadDB Storage (Native Single Layer Indirection):



Monad resolves this storage bottleneck through MonadDB, a custom state database designed specifically for storing MPT nodes on disk.<sup>30</sup> MonadDB implements a Patricia Trie natively, both on-disk and in-memory, eliminating the double-layer indirection of traditional databases.<sup>4</sup>

Furthermore, to enable concurrent parallel execution, reads must be non-blocking.<sup>30</sup> Traditional databases lack proper asynchronous I/O (async I/O) support, which forces clients to spawn large numbers of kernel threads to handle pending I/O requests.<sup>30</sup> MonadDB fully utilizes the

latest Linux kernel support for async I/O via the `io_uring` interface, allowing the database to execute concurrent reads without blocking execution threads or incurring context-switching overhead.<sup>30</sup>

To minimize overhead further, MonadDB provides operators the option to bypass the filesystem.<sup>30</sup> Modern filesystems introduce significant performance costs due to block allocation, fragmentation, read/write amplification, and metadata management.<sup>30</sup> Bypassing the filesystem and running MonadDB directly on raw block devices allows the database to fully utilize SSD hardware performance.<sup>30</sup>

MonadDB implements its own indexing system based on the Patricia trie structure, coordinating concurrent access between a single writer (execution) and multiple readers (consensus and RPC).<sup>30</sup> This native architecture reduces the number of page reads from SSDs to perform a lookup, ensuring that disk reads (costing 40-100 microseconds) can execute in parallel without stalling execution threads.<sup>4</sup>

<b>Storage Parameter</b>	<b>LevelDB / RocksDB (Traditional Clients)</b>	<b>MonadDB (Custom Radix Store)</b>
<b>Data Structure Model</b>	Embedded LSM-Tree (Nested indirection) <sup>30</sup>	Native Radix Tree (No nested indirection) <sup>4</sup>
<b>State Format Compatibility</b>	Flat key-value store, MPT layered on top <sup>30</sup>	Optimized natively for Merkle Patricia Trie nodes <sup>30</sup>
<b>Async I/O Support</b>	Limited (Spawns synchronous kernel thread pools) <sup>30</sup>	Native async I/O (Leverages Linux <code>io_uring</code> ) <sup>30</sup>
<b>OS Filesystem Dependency</b>	Mandatory (Interacts through system file calls) <sup>30</sup>	Optional (Bypasses filesystem; runs on raw block devices) <sup>30</sup>
<b>SSD Performance Access</b>	Squeezed by OS file locks, fragmentation, amplification <sup>30</sup>	Direct raw block access (Maximizes SSD performance) <sup>30</sup>
<b>Concurrent Access Model</b>	High lock contention under multi-threaded loads	Single writer (Execution), multiple readers (Consensus, RPC) <sup>30</sup>

<b>Disk Read Latency</b>	High (Cumulative latency over multiple index layers)	40–100 microseconds per native parallel read <sup>31</sup>
--------------------------	--	--

## The Monad Gas Model: DoS Prevention and Economic Alignment

To support asynchronous execution, Monad cannot use Ethereum's standard gas billing model.<sup>32</sup> In Ethereum, a transaction is billed based on the exact amount of gas consumed during execution.<sup>32</sup> In Monad's asynchronous model, block leaders must build blocks and reach consensus on transaction ordering *before* those transactions are executed.<sup>32</sup>

If Monad billed based on actual gas used, a malicious actor could submit transactions with very high gas limits that consume virtually no gas.<sup>32</sup> This would allow them to pack a block to its gas limit during the consensus phase without paying for the reserved space, creating a severe Denial of Service (DoS) vulnerability.<sup>32</sup>

To resolve this, Monad implements a unique billing rule: **users are charged for the entire gas limit specified in the transaction, regardless of actual consumption.**<sup>32</sup>

$$\text{gas\_paid} = \text{gas\_limit} \times \text{price\_per\_gas}$$

This model prevents DoS attacks by forcing users to pay for any block space they reserve.<sup>32</sup> However, this shift introduces significant friction for developers and users. If a transaction reverts during execution (e.g., due to slippage limits or a sold-out NFT), standard EVM wallets like MetaMask fallback to estimating an extremely high gas limit.<sup>32</sup> In Ethereum, the user only pays for the gas used up to the reversion point. In Monad, the user is charged for the entire, inflated gas limit, resulting in severe and unexpected financial penalties.<sup>32</sup>

Additionally, during protocol deployments, asynchronous execution delays can cause transactions to be rejected by the mempool if they use standard EVM gas estimation limits.<sup>17</sup> Developers are forced to apply massive gas estimate overrides (such as using a `--gas-estimate-multiplier 200` override) to ensure successful execution, creating unpredictable deployment costs.<sup>17</sup>

To align execution costs with physical hardware constraints, Monad has reweighted the gas pricing of cold state access opcodes.<sup>32</sup> While Monad's execution engine makes computation up to 300x cheaper, the physical cost of reading data from disk (state reads) remains bounded by SSD hardware speeds.<sup>32</sup>

If Monad kept Ethereum's original gas weights, storage would become 300x cheaper in relative terms, making the network highly vulnerable to storage-based DoS-style spam attacks.<sup>32</sup> To

prevent this, Monad has adjusted cold state access costs upward to reflect the physical cost of disk reads.<sup>32</sup> This is summarized in the opcode repricing matrix below<sup>32</sup>:

<b>Opcode / Precompile</b>	<b>Affected Instructions</b>	<b>Ethereum Gas Cost</b>	<b>Monad Gas Cost</b>	<b>Multiplier</b>
<b>Cold Account Access</b>	BALANCE, EXTCODESIZE, EXTCODECOPY, EXTCODEHASH, CALL, CALLCODE, DELEGATECALL, STATICCALL, SELFDESTRUCT	2,600	<b>10,100</b>	3.88x
<b>Cold Storage Access</b>	SLOAD, SSTORE	2,100	<b>8,100</b>	3.85x
<b>Warm Account Access</b>	BALANCE, EXTCODESIZE, etc. (Warm state cache)	100	<b>100</b>	1.0x (No change)
<b>Warm Storage Access</b>	SLOAD, SSTORE (Warm state cache)	100	<b>100</b>	1.0x (No change)
<b>ecRecover</b>	0x01 precompile (Signature recovery)	3,000	<b>6,000</b>	2.0x

<b>ecAdd</b>	0x06 precompile (Elliptic curve addition)	150	<b>300</b>	2.0x
<b>ecMul</b>	0x07 precompile (Elliptic curve multiplication)	6,000	<b>30,000</b>	5.0x
<b>ecPairing</b>	0x08 precompile (Elliptic curve pairing)	45,000	<b>225,000</b>	5.0x
<b>blake2f</b>	0x09 precompile (Cryptographic hashing)	rounds × 1	<b>rounds × 2</b>	2.0x
<b>point eval</b>	0x0a precompile (Blob transaction verification)	50,000	<b>200,000</b>	4.0x

To manage base fee adjustments at 400 ms block intervals, Monad uses a variance-aware controller inspired by adaptive optimization algorithms like RMSprop and Adam.<sup>32</sup> This controller is defined by the following mathematical equations<sup>32</sup>:

$$\text{block\_gas}_k = \sum_{\text{tx} \in \text{block}_k} \text{gas\_limit}_{\text{tx}}$$

$$\text{base\_price\_per\_gas}_{k+1} = \max \left\{ \text{min\_base\_price\_per\_gas}, \text{base\_price\_per\_gas}_k \cdot \exp \left( \eta_k \cdot \frac{t}{\text{blo}}$$

The dynamic learning rate ( $\eta_k$ ), which dampens step size adjustments under volatile demand, is calculated using the historical trend and moment of gas usage<sup>32</sup>:

$$\eta_k = \frac{\text{max\_step\_size} \cdot \epsilon}{\epsilon + \sqrt{\text{moment}_k - \text{trend}_k^2}}$$

$$\text{trend}_{k+1} = \beta \cdot \text{trend}_k + (1 - \beta) \cdot (\text{target} - \text{block\_gas}_k)$$

$$\text{moment}_{k+1} = \beta \cdot \text{moment}_k + (1 - \beta) \cdot (\text{target} - \text{block\_gas}_k)^2$$

Here, the target utilization is set to 80% (160 million gas) of the 200 million gas block limit, the smoothing factor  $\beta$  is set to 0.96 (equivalent to a 17-block half-life), and the maximum step size is capped at  $1/28$ .<sup>32</sup> This mathematical model prevents the base fee from becoming too volatile ("twitchy") during short demand spikes, but it increases fee complexity and makes it harder for developers to predict transaction costs.<sup>32</sup>

---

## State Expansion, Hardware Demands, and Node Centralization

With a testnet block gas limit of 150 million gas and 400 ms block times, Monad processes an unprecedented 375 million gas per second.<sup>4</sup> Under production workloads, this throughput will accelerate state expansion to extreme levels.<sup>1</sup> State expansion refers to the accumulation of contract code, accounts, and storage slots that must be held in a node's active memory or storage to validate transactions.<sup>33</sup>

Monad's high write throughput can generate up to 6TB of historical archive data per day.<sup>20</sup> This state expansion creates a massive operational footprint for node operators, leading to serious performance degradation and node centralization.<sup>1</sup>

To achieve sufficient state read/write speeds, Monad enforces demanding hardware requirements.<sup>35</sup> Operating a node requires a dedicated, bare-metal server; virtualization (e.g., AWS, GCP, Azure) is strictly discouraged due to the latency and context-switching overhead introduced by hypervisors.<sup>35</sup>

Furthermore, standard full nodes only retain a rolling window of approximately 40,000 blocks of state history (about 4.5 hours of operations) before aggressively pruning.<sup>17</sup> A historical archive node, which is essential for developer indexing, blockchain explorers, and historical analytics, requires over 16TB of NVMe storage and takes up to a week to bootstrap and sync.<sup>17</sup>

The physical limits of Monad's hardware profiles are detailed below:

Technical Parameter	Full Node Specification	Active Validator Node	Archive Node Specification
<b>CPU Architecture</b>	16-Core, ≥ Base (e.g., AMD Ryzen 9950x) <sup>35</sup>	16-Core, ≥ Base (e.g., AMD Ryzen 9950x) <sup>35</sup>	16-Core, ≥ Base (e.g., AMD Ryzen 9950x) <sup>35</sup>
<b>System Memory</b>	≥ RAM <sup>35</sup>	≥ RAM <sup>35</sup>	≥ RAM <sup>35</sup>
<b>Storage Configuration</b>	2TB PCIe Gen4x4 NVMe SSD for TrieDB <sup>35</sup>	2TB TrieDB NVMe SSD + 500GB MonadBFT NVMe SSD <sup>35</sup>	≥ Dedicated Enterprise NVMe SSD <sup>17</sup>
<b>Bandwidth (Ingress/Egress)</b>	≥ Sustained <sup>4</sup>	≥ Sustained <sup>4</sup>	≥ Dedicated Burstable
<b>Virtualization Compatibility</b>	Unstable (Highly discouraged) <sup>35</sup>	Incompatible (Slashes performance) <sup>35</sup>	Incompatible (Extreme hypervisor I/O bottlenecks)
<b>Consensus Stake Threshold</b>	N/A (Non-voting listener) <sup>10</sup>	≥ Total Stake (Top 200 restriction) <sup>12</sup>	N/A (Query service layer)
<b>Bootstrap / Sync Duration</b>	~	~	Up to 7 Days (Sustained I/O bottleneck) <sup>17</sup>

These demanding requirements introduce a clear decentralization ceiling.<sup>1</sup> The high cost of bare-metal hardware and dedicated bandwidth restricts node operation to professional data centers and well-capitalized institutions, which could lead to validator centralization.<sup>1</sup>

Furthermore, because Monad's active consensus validator set is limited to the top 200 nodes by total stake (with a minimum self-stake of 100,000 MON and a total stake requirement of 10 million MON)<sup>12</sup>, block production is highly concentrated. If historical data access is restricted to a small number of centralized archive providers, the network's trustlessness could be

compromised.<sup>20</sup>

---

## Nuanced Conclusions and Actionable Recommendations

The technological design of the Monad blockchain pushes physical hardware and distributed network limits to deliver massive transaction throughput.<sup>3</sup> However, this hyperscaling introduces deep systemic trade-offs, storage bottlenecks, and execution vulnerabilities.<sup>1</sup>

To optimize performance, reduce developer friction, and preserve decentralization, Category Labs should implement the following targeted engineering modifications:

### 1. Implement Proactive Conflict-Aware Scheduling

Monad's current optimistic execution model (2PE) suffers from high abort and re-execution rates under heavy state contention, dragging down performance to sequential speeds.<sup>26</sup>

To resolve this, Monad should integrate a proactive scheduling layer similar to SupraBTM.<sup>26</sup> Before executing transactions, a lightweight scheduler should analyze transaction signatures and mempool state to detect conflicts.<sup>26</sup>

If the anticipated conflict ratio exceeds a dynamically calibrated threshold, the system should adaptively schedule the conflicting subset of transactions or group them into a dependency DAG prior to execution.<sup>26</sup> This hybrid approach would prevent threads from executing transactions that are guaranteed to abort, saving valuable CPU cycles under heavy load.<sup>22</sup>

### 2. Deploy Speculative RPC State-Tracking Proxies

The loss of immediate read-your-own-writes (RYW) consistency due to deferred execution creates major integration challenges for client-side applications, trading bots, and keepers.<sup>17</sup>

To restore seamless integration, Monad's RPC nodes should implement a local, speculative state-tracking proxy.<sup>19</sup> This proxy would maintain a short-lived memory projection of pending transactions, routing read requests (such as `eth_call`) through a speculative execution sandbox that applies the client's pending transactions onto the current committed state.<sup>19</sup>

Additionally, session-based routing should be enforced to pin a client's read requests to the same node that processed their writes, guaranteeing immediate local consistency without impacting global validator performance.<sup>18</sup>

### 3. Integrate Native State Expiry and Eviction Protocols

The high write throughput of Monad's 375 million gas/second rate accelerates state expansion, creating severe storage and bandwidth barriers for node operators.<sup>1</sup>

To prevent runaway state bloat and reduce hardware requirements, Monad should build an

active state-expiry protocol directly into MonadDB.<sup>30</sup> Accounts and contract storage keys that have not been accessed within a configured epoch window (e.g., 30 days) should be evicted from the active on-disk Patricia Trie and moved to cold, decentralized storage layers.<sup>34</sup>

To access these accounts, users would submit a cryptographic proof of the historical state alongside their transaction.<sup>34</sup> This active pruning would cap the active TrieDB storage requirement to a predictable size, allowing validators to run on affordable commodity NVMe drives without losing performance.<sup>4</sup>

#### 4. Provide Post-Execution Gas Refunds

Charging users for their entire gas limit (rather than actual consumption) prevents DoS attacks under deferred execution, but it penalizes honest users when transactions revert due to standard conditions.<sup>32</sup>

To reduce this friction, Monad should implement a secure post-execution gas refund mechanism.<sup>32</sup> During the serial commit phase (Phase 2), if a transaction is proven to be honest (e.g., it did not revert due to state-exhaustion attacks and actual consumption was far below the specified limit), the protocol should refund up to 50% of the unused gas limit directly to the user's balance.<sup>25</sup>

This would protect users from severe over-billing while maintaining robust DoS protection across the network.<sup>32</sup>

#### Works cited

1. What Is Monad? The High-Throughput EVM Layer-1 With Parallel Execution - Atomic Wallet, accessed May 17, 2026, <https://atomicwallet.io/academy/articles/what-is-monad>
2. Parallel Execution & Monad, accessed May 17, 2026, <https://www.monad.xyz/announcements/parallel-execution-monad>
3. Monad Review: The fast parallel EVM - Rango Exchange, accessed May 17, 2026, <https://rango.exchange/learn/market-trends/monad-blockchain-review>
4. How Monad Works, accessed May 17, 2026, <https://blog.monad.xyz/blog/how-monad-works>
5. Monad for Developers - Monad Documentation, accessed May 17, 2026, <https://docs.monad.xyz/introduction/monad-for-developers>
6. Monad: A Supercharged EVM Layer 1 - Blockworks Research, accessed May 17, 2026, [https://app.blockworksresearch.com/unlocked/monad\\_a-supercharged-evm-layer-1](https://app.blockworksresearch.com/unlocked/monad_a-supercharged-evm-layer-1)
7. High-Performance EVM Layer-1 for Scalable DeFi - Infura, accessed May 17, 2026, <https://www.infura.io/networks/monad>
8. MonadBFT: Fast, Responsive, Fork-Resistant Streamlined Consensus - arXiv, accessed May 17, 2026, <https://arxiv.org/html/2502.20692v2>
9. MonadBFT Analysis (Part 1): Resolving Tail Fork Issues | Gate Learn, accessed May

- 17, 2026,  
<https://www.gate.com/learn/articles/monad-bft-analysis-part-1-resolving-tail-for-k-issues/8817>
10. How Monad Works - HackMD, accessed May 17, 2026,  
<https://hackmd.io/@keone/how-monad-works>
  11. MonadBFT: Fast, Responsive, Fork-Resistant Streamlined Consensus - ResearchGate, accessed May 17, 2026,  
[https://www.researchgate.net/publication/389510445\\_MonadBFT\\_Fast\\_Responsive\\_Fork-Resistant\\_Streamlined\\_Consensus](https://www.researchgate.net/publication/389510445_MonadBFT_Fast_Responsive_Fork-Resistant_Streamlined_Consensus)
  12. Monad First Look: Hyperscaling the EVM - Figment.io, accessed May 17, 2026,  
<https://www.figment.io/insights/monad-first-look-hyperscaling-the-evm/>
  13. MonadBFT: Fast, Responsive, Fork-Resistant Streamlined Consensus - arXiv, accessed May 17, 2026, <https://arxiv.org/pdf/2502.20692>
  14. Monad's Road to Mainnet | Everstake, accessed May 17, 2026,  
<https://everstake.one/resources/blog/monads-road-to-mainnet>
  15. Pioneering Parallelisation with Monad - Luganodes, accessed May 17, 2026,  
<https://luganodes.com/blog/PioneeringParallelisationMonad>
  16. Monad Deferred Execution: Unlocking Scalability and Efficiency in Blockchain Transactions, accessed May 17, 2026,  
[https://medium.com/@kintsu\\_xyz/monad-deferred-execution-unlocking-scalability-and-efficiency-in-blockchain-transactions-98afe27f38d8](https://medium.com/@kintsu_xyz/monad-deferred-execution-unlocking-scalability-and-efficiency-in-blockchain-transactions-98afe27f38d8)
  17. AL Technical Analysis Aave V3.6 <> Monad - Risk - Aave, accessed May 17, 2026,  
<https://governance.aave.com/t/al-technical-analysis-aave-v3-6-monad/24402>
  18. Read-Your-Writes Consistency - Arpit Bhayani, accessed May 17, 2026,  
<https://arpitbhayani.me/blogs/read-your-write-consistency/>
  19. A Deep Dive on Read Your Own Writes Consistency - DZone, accessed May 17, 2026, <https://dzone.com/articles/read-your-own-writes-consistency>
  20. How Monad handles 10K+ RPS when legacy infrastructure won't work - Goldsky, accessed May 17, 2026, <https://goldsky.com/case-studies/monad>
  21. Parallel Execution - Monad Documentation, accessed May 17, 2026,  
<https://docs.monad.xyz/monad-arch/execution/parallel-execution>
  22. Parallel Transaction Execution in Blockchains | Vlad Temian - Centrist Tech Optimist, accessed May 17, 2026,  
<https://blog.vtemian.com/post/parallel-transaction-execution-in-blockchains/>
  23. Sui vs Solana: Comprehensive Blockchain Comparison & Analysis - Ledger, accessed May 17, 2026,  
<https://www.ledger.com/academy/topics/blockchain/sui-vs-solana>
  24. [CryptoTimes] Blockchain & Technology Enabling Parallel Execution - DeSpread Research, accessed May 17, 2026,  
<https://research.despread.io/crypto-times-parallel-execution/>
  25. Deep Dive into Monad's Architecture - Chorus One, accessed May 17, 2026,  
<https://chorus.one/articles/deep-dive-into-monads-architecture>
  26. Supra vs. Monad: Towards the Best Parallel Execution of the ..., accessed May 17, 2026, <https://supra.com/academy/supra-vs-monad/>
  27. Monad Scanner. Monad Execution & State Risk Scanner... | by Chain Digital |

- Medium, accessed May 17, 2026,  
<https://medium.com/@chaindigital/monad-scanner-f0e4e59ba9c2>
28. Solana vs. Aptos vs. Sui: Which Layer-1 Is Best Positioned for Growth? - SorooshX, accessed May 17, 2026,  
<https://www.sorooshx.com/blog/solana-vs-aptos-vs-sui-layer-1-growth>
  29. What Is Monad: Complete Parallelized EVM Layer 1 Guide (2026) - by DEXTools, accessed May 17, 2026,  
<https://www.dextools.io/tutorials/what-is-monad-parallelized-evm-l1-guide-2026>
  30. MonadDb - Monad Documentation, accessed May 17, 2026,  
<https://docs.monad.xyz/monad-arch/execution/monaddb>
  31. Deep Dive into MonadDB. The limitations of traditional... | by MonadWhisper - Medium, accessed May 17, 2026,  
<https://medium.com/@monadwhisper/deep-dive-into-monaddb-f3aa1c6532f5>
  32. Understanding Monad's Gas Model: Smarter, Faster, and Safer | by ..., accessed May 17, 2026,  
<https://medium.com/@huseinrasti/understanding-monads-gas-model-smarter-faster-and-safer-3b2f6319ff32>
  33. How to Raise the Gas Limit, Part 1: State Growth - Paradigm, accessed May 17, 2026,  
<https://www.paradigm.xyz/2024/03/how-to-raise-the-gas-limit-1>
  34. What is State Bloat in Blockchains? - Nervos Network, accessed May 17, 2026,  
[https://www.nervos.org/knowledge-base/state\\_bloat\\_blockchain\\_\(explainCKBot\)](https://www.nervos.org/knowledge-base/state_bloat_blockchain_(explainCKBot))
  35. Hardware Requirements - Monad Documentation, accessed May 17, 2026,  
<https://docs.monad.xyz/node-ops/hardware-requirements>